

Lecture 21: Communication

UBC 2022-23

Instructor: Mathias Lécuyer

Imports

```
In [2]: 1 import matplotlib.pyplot as plt
2 import numpy as np
3 import pandas as pd
4 from sklearn.compose import ColumnTransformer, TransformedTargetRegressor
5 from sklearn.dummy import DummyClassifier, DummyRegressor
6 from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor
7 from sklearn.impute import SimpleImputer
8 from sklearn.linear_model import Ridge
9 from sklearn.metrics import log_loss
10 from sklearn.model_selection import (
11     GridSearchCV,
12     cross_val_score,
13     cross_validate,
14     train_test_split,
15 )
16 from sklearn.pipeline import Pipeline, make_pipeline
17 from sklearn.preprocessing import (
18     MinMaxScaler,
19     OneHotEncoder,
20     OrdinalEncoder,
21     StandardScaler,
22 )
23
24 plt.rcParams["font.size"] = 16
```

Learning objectives

- When communicating about applied ML, tailor an explanation to the intended audience.
- Apply best practices of technical communication, such as bottom-up explanations and reader-centric writing.
- Given an ML problem, analyze the decision being made and the objectives.
- Avoid the pitfall of thinking about ML as coding in isolation; build the habit of relating your work to the surrounding context and stakeholders.
- Interpret a confidence score or credence, e.g. what does it mean to be 5% confident that a statement is true.
- Maintain a healthy skepticism of `predict_proba` scores and their possible interpretation as credences.
- Be careful and precise when communicating confidence to stakeholders in an ML project.

- Identify misleading visualizations.

Attribution

- The first part of this lecture is adapted from [DSCI 542 \(https://github.com/UBC-MDS/DSCI_542_comm-arg\)](https://github.com/UBC-MDS/DSCI_542_comm-arg), created by [David Laing \(https://davidklaing.com/\)](https://davidklaing.com/).
- The visualization component of this lecture benefitted from discussions with [Firas Moosvi \(http://firas.moosvi.com/\)](http://firas.moosvi.com/) about his course, [DSCI 531 \(https://github.com/UBC-MDS/DSCI_531_viz-1\)](https://github.com/UBC-MDS/DSCI_531_viz-1).

Why should we care about effective communication?

- Most ML practitioners work in an organization with >1 people.
- There will very likely be stakeholders other than yourself.
- Those people need to understand what you're doing because:
 - their state of mind may change the way you do things (see below)
 - your state of mind may change the way they do things (interpreting your results)

ML suffers from some particular communication issues:

- overstating one's results / unable to articulate the limitations
- unable to explain the predictions
- Can we trust test error?
- Why did CatBoost make that prediction?
- What does it mean if `predict_proba` outputs 0.9?

These issues are there because these things are actually very hard to explain!

Activity: explaining GridSearchCV

Below are two possible explanations of `GridSearchCV` pitched to different audiences. Read them both and then follow the instructions at the end.

Explanation 1

Machine learning algorithms, like an airplane's cockpit, typically involve a bunch of knobs and switches that need to be set.



For example, check out the documentation of the popular random forest algorithm [here](https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html) (<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>). Here's a list of the function arguments, along with their default values (from the documentation):

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, criterion='gini',
max_depth=None, min_samples_split=2, min_samples_leaf=1,
min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None, bootstrap=True,
oob_score=False, n_jobs=None, random_state=None, verbose=0,
warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

Holy cow, that's a lot of knobs and switches! As a machine learning practitioner, how am I supposed to choose `n_estimators` ? Should I leave it at the default of 100? Or try 1000? What about `criterion` or `class_weight` for that matter? Should I trust the defaults?

Enter [GridSearchCV](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html) (https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html) to save the day. The general strategy here is to choose the settings that perform best on the specific task of interest. So I can't say `n_estimators=100` is better than `n_estimators=1000` without knowing what problem I'm working on. For a specific problem, you usually have a numerical score that measures performance. `GridSearchCV` is part of the popular [scikit-learn](https://scikit-learn.org/) (<https://scikit-learn.org/>) Python machine learning library. It works by searching over various settings and tells you which one worked best on your problem.

The "grid" in "grid search" comes from the fact that it tries all possible combinations on a grid. For example, if you want it to consider setting `n_estimators` to 100, 150 or 200, and you want it to consider setting `criterion` to 'gini' or 'entropy', then it will search over all 6 possible combinations in a grid of 3 possible values by 2 possible values:

`criterion='gini' criterion='entropy'`

```
In [3]: 1 # imports
2 from sklearn import datasets
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.model_selection import GridSearchCV
5
6 # load a dataset
7 data = datasets.load_digits()
8 X = data["data"]
9 y = data["target"]
10
11 # set up the grid search
12 grid_search = GridSearchCV(
13     RandomForestClassifier(random_state=123),
14     param_grid={"n_estimators": [100, 150, 200], "criterion": ["gini",
15 ]
16
17 # run the grid search
18 grid_search.fit(X, y)
19 grid_search.best_params_
```

```
Out[3]: {'criterion': 'gini', 'n_estimators': 100}
```

As we can see from the output above, the grid search selected `criterion='gini'`, `n_estimators=100`, which was one of our 6 options above (specifically Option 1).

By the way, these "knobs" we've been setting are called [hyperparameters](https://en.wikipedia.org/wiki/Hyperparameter_(machine_learning)) ([https://en.wikipedia.org/wiki/Hyperparameter_\(machine_learning\)](https://en.wikipedia.org/wiki/Hyperparameter_(machine_learning))) and the process of setting these hyperparameters automatically is called [hyperparameter optimization](https://en.wikipedia.org/wiki/Hyperparameter_optimization) (https://en.wikipedia.org/wiki/Hyperparameter_optimization) or *hyperparameter tuning*.

~400 words, not including code.

Explanation 2

<https://medium.com/datadriveninvestor/an-introduction-to-grid-search-ff57adcc0998>
(<https://medium.com/datadriveninvestor/an-introduction-to-grid-search-ff57adcc0998>)

~400 words, not including code.

Discussion questions:

- What do you like about each explanation?
- What do you dislike about each explanation?
- What do you think is the intended audience for each explanation?
- Which explanation do you think is more effective overall for someone on Day 1 of CPSC 330?
- Each explanation has an image. Which one is more effective? What are the pros/cons?
- Each explanation has some sample code. Which one is more effective? What are the pros/cons?

After you're done reading, take ~5 min to consider the discussion questions above. Paste your answer to **at least one** of the above questions in the [Google jamboard](https://jamboard.google.com/d/1WbJTtNi-qt4EjvONcyO-CUxOhGxXhaRwYto28fC5Kzs/edit?usp=sharing) (<https://jamboard.google.com/d/1WbJTtNi-qt4EjvONcyO-CUxOhGxXhaRwYto28fC5Kzs/edit?usp=sharing>) under the appropriate question heading.

Principles of good explanations

Concepts *then* labels, not the other way around

The first explanation start with an analogy for the concept (and the label is left until the very end):

Machine learning algorithms, like an airplane's cockpit, typically involve a bunch of knobs and switches that need to be set.

In the second explanation, the first sentence is wasted on anyone who doesn't already know what "hyperparameter tuning" means:

Grid search is the process of performing hyper parameter tuning in order to determine the optimal values for a given model.

The effectiveness of these different statements depend on your audience.

See [this video](https://twitter.com/ProfFeynman/status/899963856549625858?s=20) (<https://twitter.com/ProfFeynman/status/899963856549625858?s=20>): "I learned very early the difference between knowing the name of something and knowing something." - Richard Feynman.

Bottom-up explanations

The [Curse of Knowledge \(https://en.wikipedia.org/wiki/Curse_of_knowledge\)](https://en.wikipedia.org/wiki/Curse_of_knowledge) leads to *top-down* explanations:

- When you know something well, you think about things in the context of all your knowledge.
- Those lacking the context, or frame of mind, cannot easily understand.

There is another way: *bottom-up* explanations:

When you're brand new to a concept, you benefit from analogies, concrete examples and familiar patterns.

New ideas in small chunks

The first explanation has a hidden conceptual skeleton:

1. The concept of setting a bunch of values.
2. Random forest example.
3. The problem / pain point.
4. The solution.
5. How it works - high level.
6. How it works - written example.
7. How it works - code example.
8. The name of what we were discussing all this time.

Reuse your running examples

Effective explanations often use the same example throughout the text and code. This helps readers follow the line of reasoning.

Approach from all angles

When we're trying to draw mental boundaries around a concept, it's helpful to see examples on all sides of those boundaries. If we were writing a longer explanation, it might have been better to show more, e.g.

- Performance with and without hyperparameter tuning.
- Other types of hyperparameter tuning (e.g. `RandomizedSearchCV`).

When experimenting, show the results asap

The first explanation shows the output of the code, whereas the second does not. This is easy to do and makes a big difference.

Interesting to you != useful to the reader (aka it's not about you)

Here is something which was deleted from the explanation:

Some hyperparameters, like `n_estimators` are numeric. Numeric hyperparameters are like the knobs in the cockpit: you can tune them continuously. `n_estimators` is numeric. Categorical hyperparameters are like the switches in the cockpit: they can take on (two or more) distinct values. `criterion` is categorical.

It's a very elegant analogy! But is it helpful?

And furthermore, what is my hidden motivation for wanting to include it? Elegance, art, and the pursuit of higher beauty? Or *making myself look smart*? So maybe another name for this principle could be **It's not about you**.

ML and decision-making

- There is often a wide gap between what people care about and what ML can do.
- To understand what ML can do, let's think about what **decisions** will be made using ML.

Decisions involve a few key pieces

- The **decision variable**: the variable that is manipulated through the decision.
 - E.g. how much should I sell my house for? (numeric)
 - E.g. should I sell my house? (categorical)
- The decision-maker's **objectives**: the variables that the decision-maker ultimately cares about, and wishes to manipulate indirectly through the decision variable.
 - E.g. my total profit, time to sale, etc.
- The **context**: the variables that mediate the relationship between the decision variable and the objectives.
 - E.g. the housing market, cost of marketing it, my timeline, etc.

How does this inform you as an ML practitioner?

Questions you have to answer:

- Who is the decision maker?
- What are their objectives?

- What are their alternatives?
- What is their context?
- What data do I need?

Break (10 min)

- We'll take a longer break today.
- Consider taking this time to fill out the instructor/TA evaluations if you haven't already.

Evaluation link(s):

- https://canvas.ubc.ca/courses/83420/external_tools/4732
(https://canvas.ubc.ca/courses/83420/external_tools/4732)
- https://go.blueja.io/6smkkXqkVEq_u38wYKHE6Q
(https://go.blueja.io/6smkkXqkVEq_u38wYKHE6Q)
- Here is [Mike's post on these evaluations](https://www.reddit.com/r/UBC/comments/k18qj7/teaching_evaluations_the_good_the_bad_and)
(https://www.reddit.com/r/UBC/comments/k18qj7/teaching_evaluations_the_good_the_bad_and)

Confidence and predict_proba

- What does it mean to be "confident" in your results?
- When you perform analysis, you are responsible for many judgment calls.
- [Your results will be different than others'](https://fivethirtyeight.com/features/science-isnt-broken/#part1) (<https://fivethirtyeight.com/features/science-isnt-broken/#part1>).
- As you make these judgments and start to form conclusions, how can you recognize your own uncertainties about the data so that you can communicate confidently?

What does this mean for us, when we're trying to make claims about our data?

Let's imagine that the following claim is true:

Vancouver has the highest cost of living of all cities in Canada.

Now let's consider a few beliefs we could hold:

1. Vancouver has the highest cost of living of all cities in Canada. **I am 95% sure of this.**
2. Vancouver has the highest cost of living of all cities in Canada. **I am 55% sure of this.**

The part in bold is called a [credence](https://en.wikipedia.org/wiki/Credence_(statistics)) ([https://en.wikipedia.org/wiki/Credence_\(statistics\)](https://en.wikipedia.org/wiki/Credence_(statistics))). Which belief is better?

But what if it's actually Toronto that has the highest cost of living in Canada?

1. Vancouver has the highest cost of living of all cities in Canada. **I am 95% sure of this.**
2. Vancouver has the highest cost of living of all cities in Canada. **I am 55% sure of this.**

Which belief is better now?

Conclusion: We don't just want to be right. We want to be confident when we're right and hesitant when we're wrong.

What do credences mean in practical terms?

One of two things:

- **I would accept a bet at these odds.** 99% sure means, "For the chance of winning \$1, I would bet \$99 that I'm right about this." 75% sure means, "For the chance of winning \$25, I would bet \$75 that I'm right about this."
- **Long-run frequency of correctness.** 99% sure means, "For every 100 predictions I make at this level of confidence, I would expect only 1 of them to be incorrect." 75% sure means, "For every 100 predictions I make at this level of confidence, I would expect about 25 of them to be incorrect."

It's easy enough to evaluate how good we are at being right...

But if we want to evaluate *how good we are at knowing how right we are?*

We would need to keep track of not just the correctness of our predictions, but also our confidence in those predictions.

What does this have to do with applied ML?

- What if you `predict` that a credit card transaction is fraudulent?
 - We probably want `predict_proba` a lot of the time.
- What if `predict_proba` is 0.95 in that case?
 - How confident are YOU?
- What if you forecast that avocado prices will go up next week?
 - How confident are you there?
- Or what if you predict a house price to be \$800k?
 - That is not even a true/false statement.

Loss functions

When you call `fit` for `LogisticRegression` it has these same preferences: correct and confident > correct and hesitant > incorrect and hesitant > incorrect and confident.

```
In [3]: 1 from sklearn.metrics import log_loss
```

- This is a "loss" or "error" function like mean squared error, so lower values are better.
- When you call `fit` it tries to minimize this metric.

Correct and 95% confident:

```
In [4]: 1 log_loss(y_true=np.array([0]), y_pred=np.array([[0.95, 0.05]]), labels=
```

```
Out[4]: 0.05129329438755058
```

Correct and 55% confident:

```
In [5]: 1 log_loss(y_true=np.array([0]), y_pred=np.array([[0.55, 0.45]]), labels=
```

```
Out[5]: 0.5978370007556204
```

Incorrect and 55% confident:

```
In [6]: 1 log_loss(y_true=np.array([0]), y_pred=np.array([[0.45, 0.55]]), labels=
```

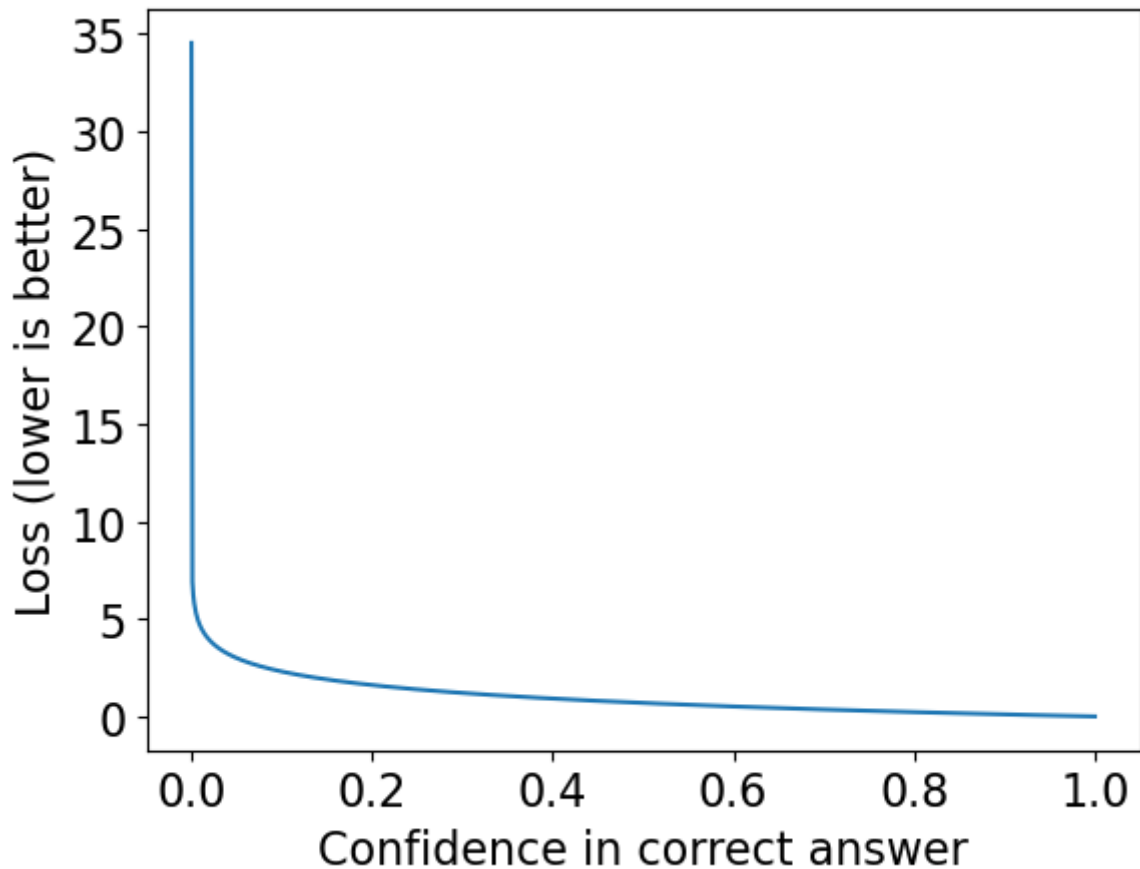
```
Out[6]: 0.7985076962177716
```

Incorrect and 95% confident:

```
In [7]: 1 log_loss(y_true=np.array([0]), y_pred=np.array([[0.05, 0.95]]), labels=
```

```
Out[7]: 2.995732273553991
```

```
In [8]: 1 grid = np.linspace(0, 1, 1000)
2 plt.plot(
3     grid,
4     [log_loss(y_true=np.array([1]), y_pred=np.array([g]), labels=(0, 1)
5 )
6 plt.xlabel("Confidence in correct answer")
7 plt.ylabel("Loss (lower is better)");
```



- Your loss goes to 0 as you approach 100% confidence in the correct answer.
- Your loss goes to infinity as you approach 100% confidence in the incorrect answer.
- (Optional) See also the very related [How to assign partial credit on an exam of true-false questions?](https://terrytao.wordpress.com/2016/06/01/how-to-assign-partial-credit-on-an-exam-of-true-false-questions/) (<https://terrytao.wordpress.com/2016/06/01/how-to-assign-partial-credit-on-an-exam-of-true-false-questions/>)

The real `LogisticRegression` is averaging this score over all training examples.

Some nice examples:

- [Scott Alexander](https://slatestarcodex.com/2019/01/22/2018-predictions-calibration-results/) (<https://slatestarcodex.com/2019/01/22/2018-predictions-calibration-results/>)
 - Look at how the plot starts at 50%. That is because being 40% confident of "X" is the same as being 60% confident of "not X".
- [Good Judgment Project](https://www.gjopen.com/) (<https://www.gjopen.com/>)

Credence Activity (time permitting: 15 min)

- Take a few minutes and assign credences or values to the claims below, in the [Google Doc](https://jamboard.google.com/d/1WbJTINi-qt4EjvONcyO-CUxOhGxXhaRwYto28fC5Kzs/edit?usp=sharing) (<https://jamboard.google.com/d/1WbJTINi-qt4EjvONcyO-CUxOhGxXhaRwYto28fC5Kzs/edit?usp=sharing>). Afterwards, we'll discuss.
 - **Do not search the answers; the point of the exercise is to evaluate how good we are at guessing.**
 - Try not to be influenced by other peoples' answers! Better to pick your answers before going to the Google Doc.
1. I am ___ % sure that the world's longest river is between 6000km and 8000km.
 2. I am ___ % sure that there is 4 to 8 liters of blood in an average adult human body.
 3. I am 99% sure that the world's tallest tree is taller than ___ m.
 4. I am 90% sure that the world's tallest tree is taller than ___ m.
 5. I am 50% sure that the world's tallest tree is taller than ___ m.

NOTE: 100% means you are completely sure the statement is true, 0% means you are completely sure the statement is false.

Answers (if you are curious):

1. Nile, 6650 km
2. 5 liters
3. (to 5) the world's tallest tree is [116.07 m](https://www.guinnessworldrecords.com/world-records/tallest-tree-living#:~:text=The%20tallest%20tree%20currently%20living,to%20try%20and%20protect%20i) (<https://www.guinnessworldrecords.com/world-records/tallest-tree-living#:~:text=The%20tallest%20tree%20currently%20living,to%20try%20and%20protect%20i>)

Visualizing your results

- Very powerful but at the same time can be misleading if not done properly.

Pre-viewing review from [Calling BS visualization videos](https://www.youtube.com/watch?v=T-5aLbNeGo0&list=PLPnZfvKID1Sje5jWxt-4CSZD7bUI4gSPS&index=30&t=0s) (<https://www.youtube.com/watch?v=T-5aLbNeGo0&list=PLPnZfvKID1Sje5jWxt-4CSZD7bUI4gSPS&index=30&t=0s>):

- Dataviz in the popular media.
 - e.g. [modern NYT](https://youtu.be/T-5aLbNeGo0?t=367) (<https://youtu.be/T-5aLbNeGo0?t=367>)
- Misleading axes.
 - e.g. [vaccines](https://youtu.be/9pNWVMxaFuM?t=299) (<https://youtu.be/9pNWVMxaFuM?t=299>)
- Manipulating bin sizes.
 - e.g. [tax dollars](https://youtu.be/zAg1wsYfwsM?t=196) (<https://youtu.be/zAg1wsYfwsM?t=196>)

- Dataviz ducks.
 - e.g. [drinking water \(https://youtu.be/rmii1hfP6d4?t=169\)](https://youtu.be/rmii1hfP6d4?t=169)
 - "look how clever we are about design" -> making it about me instead of about you (see last class)
- Glass slippers.
 - e.g. [internet marketing tree \(https://youtu.be/59teS0SUHtI?t=285\)](https://youtu.be/59teS0SUHtI?t=285)
- The principle of proportional ink.
 - e.g. [most read books \(https://youtu.be/oNhusd3xFC4?t=147\)](https://youtu.be/oNhusd3xFC4?t=147)

- [Demo of cleaning up a plot \(https://www.darkhorseanalytics.com/blog/data-looks-better-naked/\)](https://www.darkhorseanalytics.com/blog/data-looks-better-naked/)
- [Principle of proportional ink \(https://serialmentor.com/dataviz/proportional-ink.html\)](https://serialmentor.com/dataviz/proportional-ink.html) from a viz textbook.

Dataset

We'll be using [Kaggle House Prices dataset \(https://www.kaggle.com/c/home-data-for-ml-course/\)](https://www.kaggle.com/c/home-data-for-ml-course/), which we used in lecture 10. As usual, to run this notebook you'll need to download the data. For this dataset, train and test have already been separated. We'll be working with the train portion.

```
In [9]: 1 df = pd.read_csv("../data/housing-kaggle/train.csv")
        2 train_df, test_df = train_test_split(df, test_size=0.10, random_state=1)
        3 train_df.head()
```

```
Out[9]:
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	L
302	303	20	RL	118.0	13704	Pave	NaN	IR1		Lvl
767	768	50	RL	75.0	12508	Pave	NaN	IR1		Lvl
429	430	20	RL	130.0	11457	Pave	NaN	IR1		Lvl
1139	1140	30	RL	98.0	8731	Pave	NaN	IR1		Lvl
558	559	60	RL	57.0	21872	Pave	NaN	IR2		HLS

5 rows × 81 columns

```
In [10]: 1 train_df.shape
```

```
Out[10]: (1314, 81)
```

Let's separate x and y

```
In [11]: 1 x_train = train_df.drop(columns=["SalePrice"])
          2 y_train = train_df["SalePrice"]
          3
          4 x_test = test_df.drop(columns=["SalePrice"])
          5 y_test = test_df["SalePrice"]
```

Feature types

- We have mixed feature types and a bunch of missing values.
- Now, let's identify feature types and transformations.

```
In [12]: 1 drop_features = ["Id"]
2 numeric_features = [
3     "BedroomAbvGr",
4     "KitchenAbvGr",
5     "LotFrontage",
6     "LotArea",
7     "OverallQual",
8     "OverallCond",
9     "YearBuilt",
10    "YearRemodAdd",
11    "MasVnrArea",
12    "BsmtFinSF1",
13    "BsmtFinSF2",
14    "BsmtUnfSF",
15    "TotalBsmtSF",
16    "1stFlrSF",
17    "2ndFlrSF",
18    "LowQualFinSF",
19    "GrLivArea",
20    "BsmtFullBath",
21    "BsmtHalfBath",
22    "FullBath",
23    "HalfBath",
24    "TotRmsAbvGrd",
25    "Fireplaces",
26    "GarageYrBlt",
27    "GarageCars",
28    "GarageArea",
29    "WoodDeckSF",
30    "OpenPorchSF",
31    "EnclosedPorch",
32    "3SsnPorch",
33    "ScreenPorch",
34    "PoolArea",
35    "MiscVal",
36    "YrSold",
37 ]
```

```

In [13]: 1 ordinal_features_reg = [
2         "ExterQual",
3         "ExterCond",
4         "BsmtQual",
5         "BsmtCond",
6         "HeatingQC",
7         "KitchenQual",
8         "FireplaceQu",
9         "GarageQual",
10        "GarageCond",
11        "PoolQC",
12    ]
13    ordering = [
14        "Po",
15        "Fa",
16        "TA",
17        "Gd",
18        "Ex",
19    ] # if N/A it will just impute something, per below
20    ordering_ordinal_reg = [ordering] * len(ordinal_features_reg)
21    ordering_ordinal_reg

```

```

Out[13]: [['Po', 'Fa', 'TA', 'Gd', 'Ex'],
['Po', 'Fa', 'TA', 'Gd', 'Ex'],
['Po', 'Fa', 'TA', 'Gd', 'Ex'],
['Po', 'Fa', 'TA', 'Gd', 'Ex'],
['Po', 'Fa', 'TA', 'Gd', 'Ex'],
['Po', 'Fa', 'TA', 'Gd', 'Ex'],
['Po', 'Fa', 'TA', 'Gd', 'Ex'],
['Po', 'Fa', 'TA', 'Gd', 'Ex'],
['Po', 'Fa', 'TA', 'Gd', 'Ex'],
['Po', 'Fa', 'TA', 'Gd', 'Ex']]

```

```

In [14]: 1 ordinal_features_oth = [
2         "BsmtExposure",
3         "BsmtFinType1",
4         "BsmtFinType2",
5         "Functional",
6         "Fence",
7     ]
8     ordering_ordinal_oth = [
9         ["NA", "No", "Mn", "Av", "Gd"],
10        ["NA", "Unf", "LwQ", "Rec", "BLQ", "ALQ", "GLQ"],
11        ["NA", "Unf", "LwQ", "Rec", "BLQ", "ALQ", "GLQ"],
12        ["Sal", "Sev", "Maj2", "Maj1", "Mod", "Min2", "Min1", "Typ"],
13        ["NA", "MnWw", "GdWo", "MnPrv", "GdPrv"],
14    ]

```

The remaining features are categorical features.


```
In [15]: 1 categorical_features = list(  
2         set(X_train.columns)  
3         - set(numeric_features)  
4         - set(ordinal_features_reg)  
5         - set(ordinal_features_oth)  
6         - set(drop_features)  
7     )  
8 categorical_features
```

```
Out[15]: ['SaleType',  
          'MSSubClass',  
          'LandContour',  
          'CentralAir',  
          'PavedDrive',  
          'LotShape',  
          'MSZoning',  
          'MiscFeature',  
          'Alley',  
          'LotConfig',  
          'Utilities',  
          'MoSold',  
          'MasVnrType',  
          'Condition2',  
          'RoofStyle',  
          'RoofMatl',  
          'HouseStyle',  
          'Heating',  
          'GarageFinish',  
          'BldgType',  
          'SaleCondition',  
          'Foundation',  
          'Condition1',  
          'LandSlope',  
          'Neighborhood',  
          'Electrical',  
          'Exterior2nd',  
          'Exterior1st',  
          'GarageType',  
          'Street']
```

Applying feature transformations

- Since we have mixed feature types, let's use `ColumnTransformer` to apply different transformations on different features types.

```
In [16]: 1 from sklearn.compose import ColumnTransformer, make_column_transformer
2
3 numeric_transformer = make_pipeline(SimpleImputer(strategy="median"), S
4 ordinal_transformer_reg = make_pipeline(
5     SimpleImputer(strategy="most_frequent"),
6     OrdinalEncoder(categories=ordering_ordinal_reg),
7 )
8
9 ordinal_transformer_oth = make_pipeline(
10     SimpleImputer(strategy="most_frequent"),
11     OrdinalEncoder(categories=ordering_ordinal_oth),
12 )
13
14 categorical_transformer = make_pipeline(
15     SimpleImputer(strategy="constant", fill_value="missing"),
16     OneHotEncoder(handle_unknown="ignore", sparse=False),
17 )
18
19 preprocessor = make_column_transformer(
20     ("drop", drop_features),
21     (numeric_transformer, numeric_features),
22     (ordinal_transformer_reg, ordinal_features_reg),
23     (ordinal_transformer_oth, ordinal_features_oth),
24     (categorical_transformer, categorical_features),
25 )
```

Examining the preprocessed data

```
In [17]: 1 preprocessor.fit(X_train)
          2 # Calling fit to examine all the transformers.
```

```
Out[17]: ColumnTransformer(transformers=[('drop', 'drop', ['Id']),
                                         ('pipeline-1',
                                          Pipeline(steps=[('simpleimputer',
                                                             SimpleImputer(strategy
='median'))],
                                                         ('standardscaler',
                                                          StandardScaler()))],
                                         ['BedroomAbvGr', 'KitchenAbvGr', 'LotFro
ntage',
                                         'LotArea', 'OverallQual', 'OverallCon
d',
                                         'YearBuilt', 'YearRemodAdd', 'MasVnrAre
a',
                                         'BsmtFinSF1', 'BsmtFinSF2', 'BsmtUnfs
F',
                                         'TotalBsmtSF', '...
r',
                                         ['SaleType', 'MSSubClass', 'LandContou
r',
                                         'CentralAir', 'PavedDrive', 'LotShape',
                                         'MSZoning', 'MiscFeature', 'Alley',
                                         'LotConfig', 'Utilities', 'MoSold',
                                         'MasVnrType', 'Condition2', 'RoofStyl
e',
                                         'RoofMatl', 'HouseStyle', 'Heating',
                                         'GarageFinish', 'BldgType', 'SaleCondit
ion',
                                         'Foundation', 'Condition1', 'LandSlop
e',
                                         'Neighborhood', 'Electrical', 'Exterior
2nd',
                                         'Exterior1st', 'GarageType', 'Stree
t'])])])
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [19]: 1 ohe_columns = list(
2         preprocessor.named_transformers_["pipeline-4"]
3         .named_steps["onehotencoder"]
4         .get_feature_names_out(categorical_features)
5     )
6     new_columns = (
7         numeric_features + ordinal_features_reg + ordinal_features_oth + oh
8     )
```

```
In [20]: 1 X_train_enc = pd.DataFrame(
2         preprocessor.transform(X_train), index=X_train.index, columns=new_c
3     )
4     X_train_enc.head()
```

```
Out[20]:
```

	BedroomAbvGr	KitchenAbvGr	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	Yr
302	0.154795	-0.222647	2.312501	0.381428	0.663680	-0.512408	0.993969	
767	1.372763	-0.222647	0.260890	0.248457	-0.054669	1.285467	-1.026793	
429	0.154795	-0.222647	2.885044	0.131607	-0.054669	-0.512408	0.563314	
1139	0.154795	-0.222647	1.358264	-0.171468	-0.773017	-0.512408	-1.689338	
558	0.154795	-0.222647	-0.597924	1.289541	0.663680	-0.512408	0.828332	

5 rows × 263 columns

```
In [21]: 1 X_test_enc = pd.DataFrame(
2         preprocessor.transform(X_test), index=X_test.index, columns=new_col
3     )
4     X_test_enc.head()
```

```
Out[21]:
```

	BedroomAbvGr	KitchenAbvGr	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	Yr
147	0.154795	-0.222647	-0.025381	-0.085415	0.663680	-0.512408	0.993969	
676	1.372763	4.348569	-0.454788	-0.074853	-1.491366	-3.209221	-2.351883	
1304	0.154795	-0.222647	-1.790721	-0.768279	0.663680	-0.512408	1.093350	
1372	0.154795	-0.222647	0.260890	-0.058176	0.663680	0.386530	0.894587	
1427	0.154795	-0.222647	-0.454788	0.073016	-0.773017	0.386530	-0.861157	

5 rows × 263 columns

```
In [22]: 1 X_train.shape, X_test.shape
```

```
Out[22]: ((1314, 80), (146, 80))
```

Training random forests and gradient boosted trees

```
In [23]: 1 from sklearn.ensemble import GradientBoostingRegressor
```

Let's compare sklearn's GradientBoostingRegressor to RandomForestRegressor for different values of `n_estimators`.

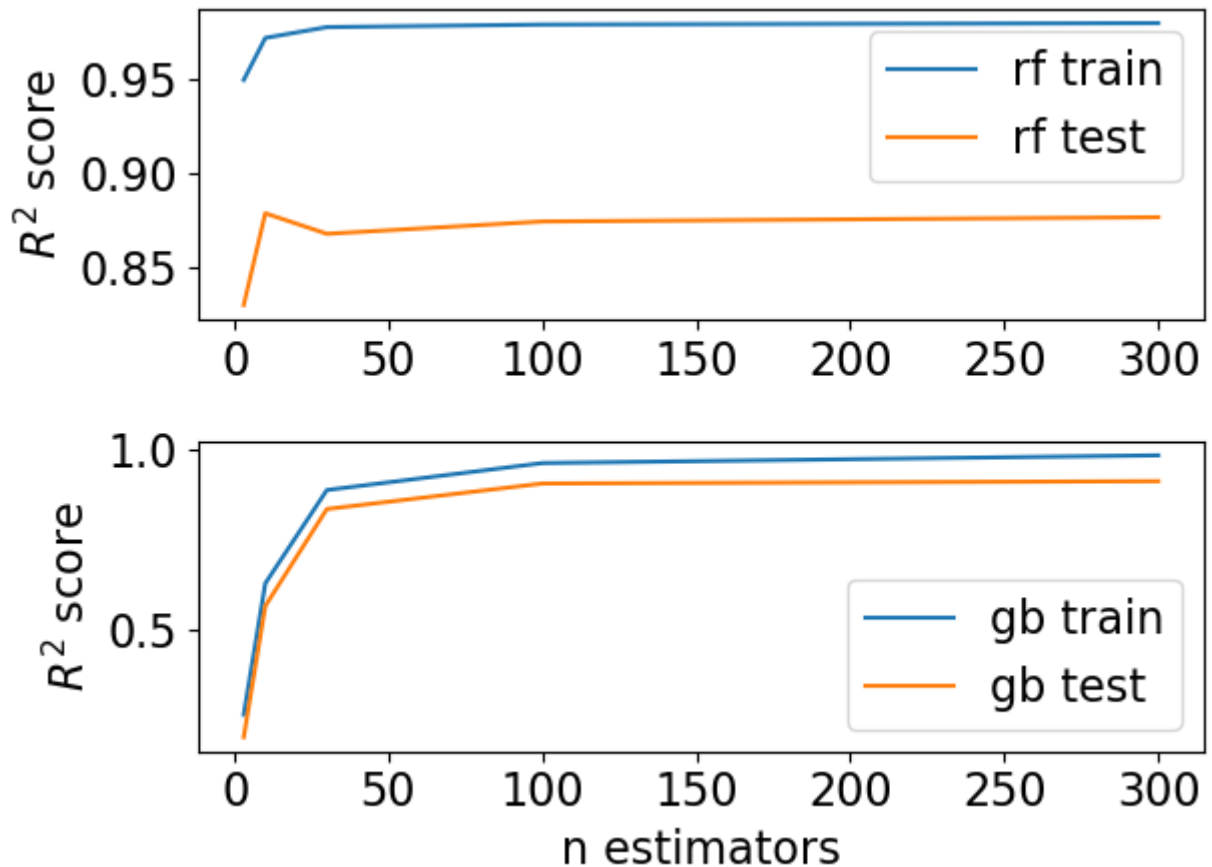
```
In [24]: 1 n_estimators_values = [3, 10, 30, 100, 300]
```

```
In [25]: 1 score_rf_train = list()
2 score_rf_test = list()
3 score_gb_train = list()
4 score_gb_test = list()
5
6 for n_estimators in n_estimators_values:
7     print(n_estimators)
8     rf = TransformedTargetRegressor(
9         RandomForestRegressor(n_estimators=n_estimators, random_state=1
10         func=np.log1p,
11         inverse_func=np.expml,
12     )
13     rf.fit(X_train_enc, y_train)
14     score_rf_train.append(rf.score(X_train_enc, y_train))
15     score_rf_test.append(rf.score(X_test_enc, y_test))
16
17     gb = TransformedTargetRegressor(
18         GradientBoostingRegressor(n_estimators=n_estimators, random_sta
19         func=np.log1p,
20         inverse_func=np.expml,
21     )
22     gb.fit(X_train_enc, y_train)
23     score_gb_train.append(gb.score(X_train_enc, y_train))
24     score_gb_test.append(gb.score(X_test_enc, y_test))
```

```
3
10
30
100
300
```

Here is a low-quality plot that is confusing and perhaps downright misleading:

```
In [26]: 1 plt.subplot(2, 1, 1)
2 plt.plot(n_estimators_values, score_rf_train, label="rf train")
3 plt.plot(n_estimators_values, score_rf_test, label="rf test")
4 plt.ylabel("$R^2$ score")
5 plt.legend()
6 plt.subplot(2, 1, 2)
7 plt.plot(n_estimators_values, score_gb_train, label="gb train")
8 plt.plot(n_estimators_values, score_gb_test, label="gb test")
9 plt.xlabel("n estimators")
10 plt.ylabel("$R^2$ score")
11 plt.legend()
12 plt.tight_layout();
```



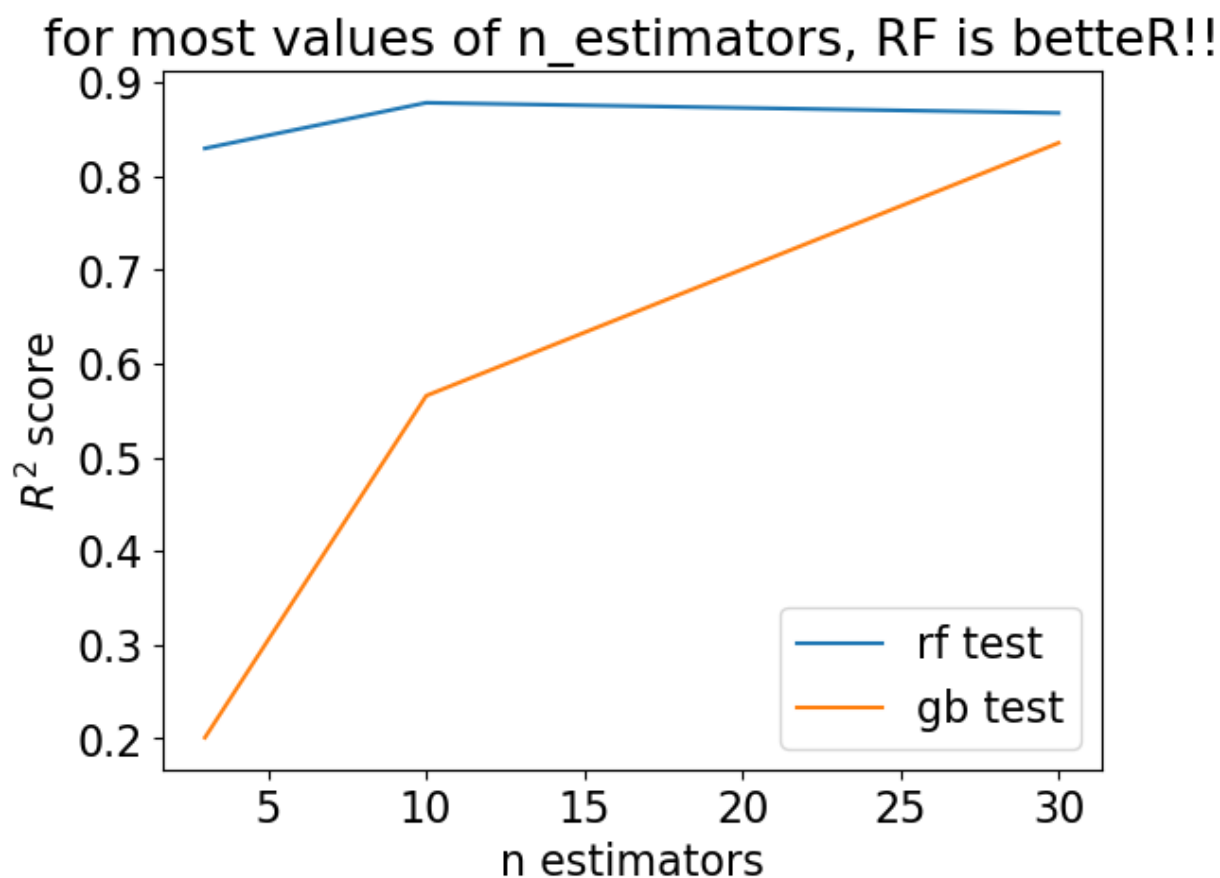
Let's create some visualizations.

- Create a visualization that makes RF look better than GB.
- Create a visualization that makes GB look better than RF.
- Create a visualization that makes RF and GB look equally good.

Here are some misleading plots.

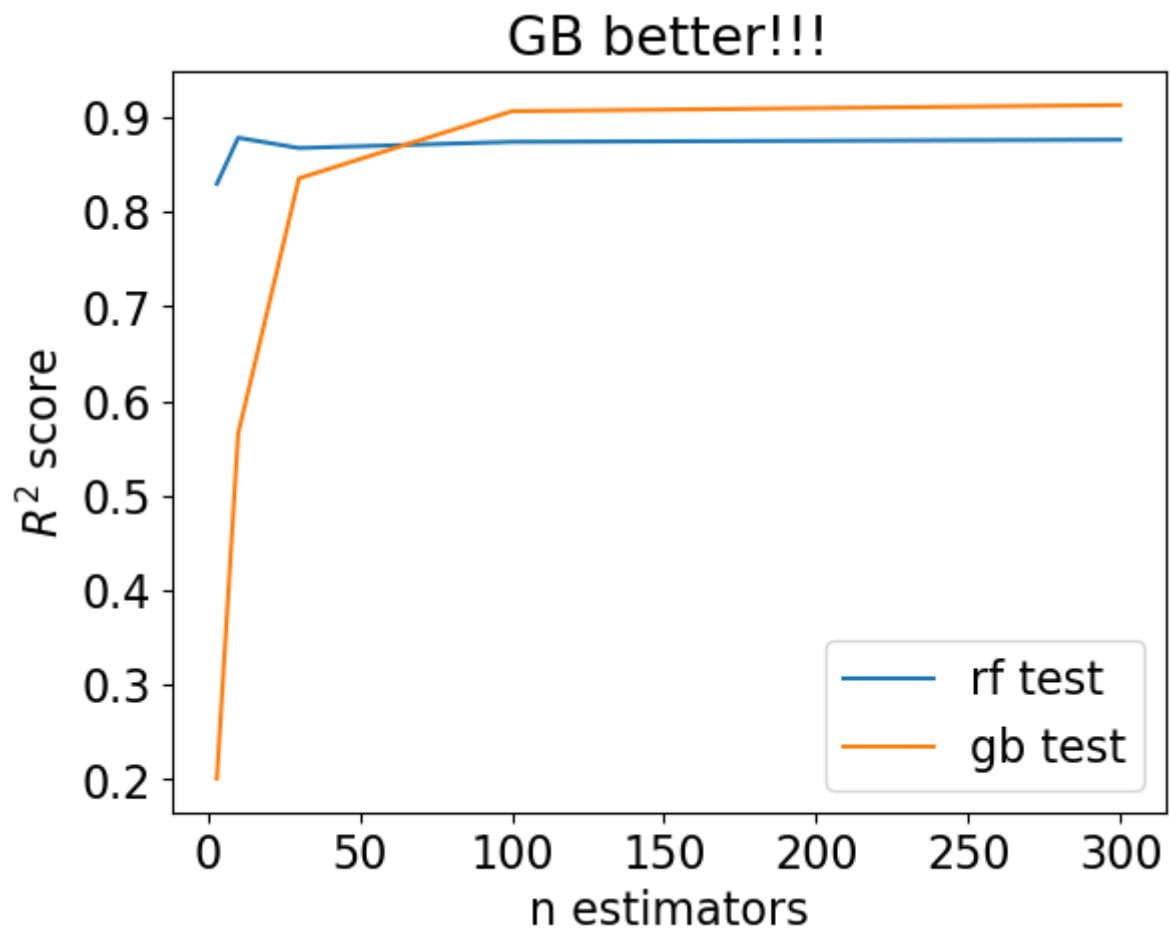
RF better than GB

```
In [27]: 1 nmax = 3
2
3 # plt.plot(n_estimators_values[:nmax], score_rf_train[:nmax], label="rf train")
4 plt.plot(n_estimators_values[:nmax], score_rf_test[:nmax], label="rf test")
5 # plt.plot(n_estimators_values[:nmax], score_gb_train[:nmax], label="gb train")
6 plt.plot(n_estimators_values[:nmax], score_gb_test[:nmax], label="gb test")
7 plt.xlabel("n estimators")
8 plt.ylabel("$R^2$ score")
9 plt.legend()
10 plt.title("for most values of n_estimators, RF is better!!");
```



GB better than RF

```
In [28]: 1 # plt.plot(n_estimators_values, score_rf_train, label="rf train")
2 plt.plot(n_estimators_values, score_rf_test, label="rf test")
3 # plt.ylabel("$R^2$ score");
4 # plt.legend();
5 # plt.subplot(2,1,2)
6 # plt.plot(n_estimators_values, score_gb_train, label="gb train")
7 plt.plot(n_estimators_values, score_gb_test, label="gb test")
8 plt.xlabel("n estimators")
9 plt.ylabel("$R^2$ score")
10 plt.legend()
11 plt.title("GB better!!!");
```



Equally good

```
In [29]: 1 nmax = 2
2
3 # plt.plot(n_estimators_values, score_rf_train, label="rf train")
4 plt.plot(n_estimators_values[:nmax], score_rf_test[:nmax], "b", label="
5 plt.ylabel("RF  $R^2$  score")
6 plt.legend(loc="lower left")
7 plt.ylim((0.8, 0.9))
8 plt.twinx()
9 # plt.plot(n_estimators_values, score_gb_train, label="gb train")
10 plt.plot(n_estimators_values[:nmax], score_gb_test[:nmax], "--r", label
11 plt.xlabel("n estimators")
12 plt.ylabel("GB  $R^2$  score")
13 plt.legend()
14 plt.ylim((-0.01, 0.70))
15 plt.title("Both equally good!!!");
```



Be critical of your visualizations and try to make them as honest as possible.

What did we learn today?

Principles of effective communication

- Concepts then labels, not the other way around.
- Bottom-up explanations.
- New ideas in small chunks.
- Reuse your running examples.
- Approaches from all angles.
- When experimenting, show the results asap.
- **It's not about you.**

- Decision variables, objectives, and context.
- How does ML fit in?
- Expressing your confidence about the results
- Misleading visualizations.

In []:

1